

**Method for creating and accessing a menu for audio content  
without using a display**

This invention relates to an audio management system that  
5 allows a user to browse through stored audio files in a very  
natural way. The invention concerns large-capacity digital  
storage-playback systems for audio content like MPEG audio  
layer 3 (MP3) players.

10

Background

Driven by the recent advances in the technologies of digital  
storage and audio compression, the problem of managing very  
15 big collections of audio files becomes predominant. For  
instance, the current generation of MP3 players contains a 10  
GB hard disk drive which enables users to store e.g. more than  
300 hours of MP3PRO music, meaning more than 4.000 titles.

20 Reliable tools are required to make those collections  
accessible to the users.

The classical way of indexing audio files is based on textual  
meta-information like title, artist, album or genre, like e.g.  
25 ID3 tags for MP3 audio files.

There are some drawbacks with this kind of organization:

1. The metadata are textual and not audio, and therefore  
cannot give a precise representation of an audio content,  
like a representative extract of the content can do.
- 30 2. Organization sorted by genre or by artist allows users to  
locate a particular piece of music. This presupposes that  
users have well-defined goals, knowing exactly what they

want to hear. The users searching strategy must be goal-driven and deterministic.

3. There are a lot of genres: for instance, the music archive mp3.com currently lists its titles under 180 different sub-genres, organized in 16 main genres. It is difficult for a user to navigate in such organization.

4. Genres are sometimes subjective because they are established a priori and not deduced from the content itself. Sometimes they are difficult to interpret.

5. A classification by genres is not able to satisfy very simple user needs like for instance "This piece of music is relaxing me. I would like to hear more like this".

The present invention is directed to overcoming these drawbacks.

#### Invention

The present invention deals with a process and system for navigating through a large amount of audio files, e.g. MP3 files, using brief representatives of the audio content. Before a user selects a music track, he can benefit from hearing a brief representative excerpt, in the following referred to as "audio thumbnail". An audio thumbnail is of sufficient length to recognize the music, e.g. 5 or 6 seconds.

The stored audio files are preprocessed in order to extract some relevant and objective descriptors. According to the invention, these descriptors are used to cluster the music tracks into perceptually homogeneous groups. From each cluster a relevant track is selected automatically or manually, or semi-automatically, and from said selected track an audio

thumbnail is extracted. Then these audio thumbnails being key phrases are arranged in a tree data structure, or table of contents, that enables the user to navigate without any visual navigation means, like display.

5

Furthermore, the audio thumbnails allow the user to navigate perceptually through the audio database, without having to remember textual elements, like title or artist names. It is particularly suited to enable users without precise idea of  
10 what they want to hear to browse their database, and to select perceptually from clusters of songs. Perceptually means here that the thumbnails address to the perception, and not memory, of users. Also, said clusters are perceptive, meaning that the structuring of the clusters is relevant for users, and thus  
15 said structuring meets real user needs.

Using this invention, users can create play lists beyond the classical music categories like pop or country.

20

#### Brief description of the drawings

Exemplary embodiments of the invention are described with reference to the accompanying drawings, which show in

25

Fig.1 an exemplary architecture of an audio reproduction system using an audio menu;

Fig.2 an exemplary user interface without display.

30

Detailed description of the invention

The present invention describes a method for creating, organizing and using audio representations for audio content.

5 The structure of the invention is shown in Figure 1. Audio tracks, usually music, are stored in a storage means S. The tracks are classified in a classifier CL and associated to a cluster of tracks C1,C2,C3. For each cluster a representative example is selected in a representative selector R. Further,  
10 an extractor X extracts a characteristic sample, or thumbnail, from said example, and the thumbnail is associated to a table of contents T. The user uses an interface I to select a first cluster represented by a first thumbnail, listen to the selected thumbnail, and decide whether to select another  
15 cluster, or select said first cluster related to said first thumbnail and then select a track belonging to said cluster, which track is then read from the storage means S and reproduced.

20 Advantageously, this approach is more perception based than previous methods, and therefore more convenient to the user. An audio-based indexing system according to the invention combines two methods that are known from other content search systems, namely the 'table-of-contents' method and the 'radio-  
25 like navigation' method.

The 'table-of-contents' method relates to a books table of contents, where short representative sequences summing up the actual text are grouped according to the structure of the  
30 books contents. This usually correlates with a logical classification into topics. Using this method for audio content means extracting parameters, or descriptors, from the audio file, following objective criteria defined below, and

then group together homogeneous tracks in clusters. From a user's point of view, these clusters make sense because their content-based nature is going farther than the a priori classification in genres. E.g. all the fragments of guitar music, coming from all genres, may be grouped together in a cluster. All the relaxing music can constitute another cluster. According to the invention, the different clusters constitute the "table of contents" of the database. And, like in a book's table of contents, there may be different levels of details, like e.g. chapter 1, chapter 1.1, etc. Like the reader can navigate from chapter to chapter, and may decide to read a chapter more in detail, the listener can navigate from cluster to cluster, or may decide to listen to more, similar music from a cluster.

The 'radio-like navigation' method relates to typical user behavior when listening to the radio. Content browsing in this context is e.g. the user scanning the FM band on a car radio, and listening to a track or switching to the next station. The invention uses this concept, with a radio station corresponding to a cluster of tracks. Then 'switch to another station' corresponds to 'select another cluster', and 'listen to a track' corresponds to 'listen to this track or to a similar track from the same cluster'.

In the following, the afore mentioned steps in creating and organizing audio representations are described in detail, the steps being performed when a track is newly added to the database, or when the database is reorganized.

In a first step descriptors are extracted from the audio track. Three types of descriptors are used, trying to be objective and still relevant for the user.

The first type of descriptors is low-level descriptors, or physical features, as being typical for signal processing methods. Examples are spectral centroid, short-time energy or  
5 short-time average zero-crossing.

The second type of descriptors is medium-level descriptors, or perceptual features, as being typically used by a musician. Examples are rhythm, e.g. binary or ternary rhythm, tonality,  
10 the kind of formation, e.g. voice or special instruments.

The third type of descriptors is high-level descriptors, or psychological and social features of the track, as being normal for the average user. Trying to minimize the subjective  
15 nature of these features, it is e.g. possible to classify music as being happy or anxious, calm or energetic. These characteristics can be assigned to a certain degree, or with a certain probability, to a piece of music, when e.g. descriptors of the previously described types are used. Also,  
20 a song can be highly memorable, can convey a certain mood or emotion, can remind the user of something, etc. This may be done automatically using supervised algorithms, i.e. with algorithms that require user interaction.

25 The second step consists of clustering the music tracks. Using the descriptors defined in the first step, the tracks can be classified into homogeneous classes. These classes are more valuable to a user than classifying music by artist or title. Unsupervised algorithms may be used to cluster the tracks into  
30 packets with similar properties. Examples of such algorithms are the K-means or Self Organizing Maps. A new cluster may be automatically generated when the dissimilarity of a newly added track, compared to existing clusters, reaches a certain

minimum level, and in that case the newly added track will be associated with the new cluster.

At this point, the tracks are classified and therefore it is possible to create a table of contents. There is no sharp classification required, e.g. it is possible to have the same track in any number of clusters. For example, one cluster may be for guitar music, while another cluster may be for calm music, and a track matching both characteristics may be associated with both clusters. In this case, both clusters may contain a link to said audio track, but the track itself needs to be stored only once.

The third step consists of automatically selecting a representative track for each cluster. Advantageously, the most representative track for a cluster is selected, using classical medoid selection. A medoid is that object of a cluster whose average dissimilarity to all objects of the cluster is minimal. Said dissimilarity can e.g. be determined using the descriptors that were extracted during the first step.

In the fourth step an audio thumbnail is created and stored for the medoid track. In another embodiment of the invention an audio thumbnail may be created and stored also for other tracks. For thumbnail creation it is evaluated which criteria are the best to characterize an audio track by a short audio sequence, the audio sequence being long enough to recognize the track, e.g. 5 or 6 seconds. In one embodiment of the invention the length of thumbnails is constant, in a second embodiment the length of thumbnails can be modified, and in a third embodiment the length of thumbnails can vary from track to track, according to the tracks descriptors. Further, in one

embodiment of the invention a thumbnail is an original sample from the track, or in another embodiment it is automatically synthesized from said track.

- 5 In the fifth step the audio thumbnails are listed in a virtual table, which can be scanned through by the user, like scanning through different radio stations. The table may be organized such that within a cluster the most relevant track, or medoid, will be found first when scanning through the table. Other
- 10 tracks within a cluster may be sorted, e.g. according to relevance. Advantageously, no graphical or textual display is required for scanning through the table of contents. The structure of the table of contents may be as follows:

15 <table of content>  
    <cluster 1>  
        <key phrase for the most relevant song (medoid)>  
        <key phrase for secondary song>  
        <key phrase 3>  
20 <cluster 2>  
    ...  
    </table of content>

- A user may decide to listen to the current track, or to
- 25 another track belonging to the same cluster and therefore being similar to said current track. Alternatively the user may decide to listen to a track from another cluster. Advantageously, only one button, or other means of command input, is required to operate the navigation system, namely
- 30 for 'Switch Cluster'. More comfortable to the user is a device with three buttons, as shown in Figure 2. One button SD is for 'Switch to a Near Cluster', one button SU is for 'Switch to a Distant Cluster', and one button M is for 'Switch to another



track from the Current Cluster'. Alternatively, it is sufficient to have only one button, if the button has more than one function, or other means of user input. Other functions controlled by user input could be e.g. random track  
5 selection or random cluster selection mode. Another function could be to successively reproduce the representatives of all clusters until the user selects one cluster, said function being advantageous because the user needs not manually scan through the table of contents.

10

Further embodiments are described in the following.

In one embodiment of the invention an audio track belongs to only one cluster, while in another embodiment an audio track  
15 may belong to more than one cluster, when the respective class criteria are not excluding each other.

In one embodiment of the invention the table of contents has only one level of clustering, like in the previously described  
20 example, while in another embodiment the table of contents can have more hierarchical levels of clusters.

In one embodiment of the invention the classification rules for audio tracks are final, while in another embodiment said  
25 rules may be modified. Said modification may happen either by an update, e.g. via internet, or by any form of user interaction, e.g. upload to PC, edit and download from PC, or by statistical or self learning methods as used e.g. by artificial intelligence. This may be implemented such that an  
30 automatic or semi-automatic reclassification with modified or enhanced rules may be performed when e.g. the number of tracks associated with one cluster is much higher than the number of tracks associated with any other cluster.

In one embodiment of the invention thumbnails may be created only for tracks representing a cluster. In another embodiment of the invention thumbnails may be created also for other  
5 tracks, e.g. tracks that fulfill a certain condition like being selected very often or very rarely, or being very long. In a third embodiment thumbnails are created for all tracks.

In one embodiment of the invention the tracks within a cluster  
10 may have a constant order, so that the user can learn after a while when a certain track comes. The order can follow the tracks relevance, or any other parameter, e.g. storage time, or frequency of selection. In another embodiment of the invention the tracks within a cluster may be unordered, or  
15 appear randomly when the user selects a cluster.

In one embodiment of the invention there is a representative track selected for each cluster, while in another embodiment it may be useful to have no representative track for one of  
20 said clusters, e.g. a cluster for favorites or a cluster for tracks not being classifiable by the employed methods.

Advantageously the described method for perception based classification and retrieval of audio contents can be used in  
25 devices, preferably portable devices, for storage and reproduction of music or other audio data, e.g. MP3 players.